

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2002-207687

(43)Date of publication of application : 26.07.2002

(51)Int.Cl. G06F 13/14
G06F 15/177

(21)Application number : 2000-384374 (71)Applicant : INTERNATL BUSINESS MACH
CORP <IBM>

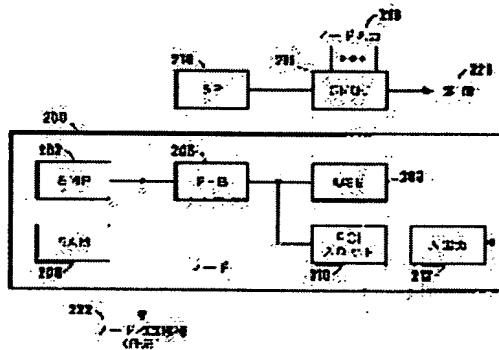
(22)Date of filing : 18.12.2000 (72)Inventor : BEALKOWSKI RICHARD
PATRICK M BRAND

(54) SYSTEM AND METHOD FOR DYNAMIC INPUT/OUTPUT ALLOCATION IN SORTED COMPUTER SYSTEM

(57)Abstract:

PROBLEM TO BE SOLVED: To provide a system and method for allowing a plurality of nodes in a multiprocessor system to share a set of input/output device.

SOLUTION: The system/device is provided with a cabinet input/output controller (CI/OC 216) which controls communication between a multiprocessor system nodes and a common input/output device 212 so as to allow the individual nodes to exclusively access one or a plurality of its target device. The respective nodes communicate with the CI/OC 216 through a service processor and the CI/OC 216 connects the various input/output devices and the USB controller 208 of the nodes with each other.



LEGAL STATUS

[Date of request for examination] 18.12.2000

[Date of sending the examiner's decision of rejection] 02.03.2004

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

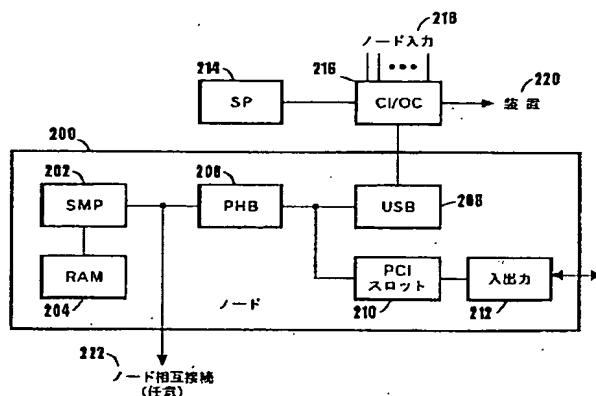
[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's



【特許請求の範囲】

【請求項1】それぞれが少なくとも1つのシステム・プロセッサおよび関連メモリを有し、それぞれがそれぞれの装置ポートを介して通信するように接続された、複数のシステム・ポートと、
前記装置ポートの少なくとも1つに接続された、装置コントローラと、

前記装置コントローラに接続された、少なくとも1つの周辺装置とを含み、前記コンピュータ・システムが、直接接続を介してシステム・ノードから前記装置コントローラへ要求を送るステップと、

前記装置コントローラによって、前記システム・ノードと前記周辺装置との間の排他的接続を確立するステップと、

前記システム・ノードによって前記周辺装置を操作するステップとを実行するように構成されるコンピュータ・システム。

【請求項2】前記システム・ノードが、対称マルチプロセッサ・システムとして動作する、請求項1に記載のシステム。

【請求項3】前記システム・ノードのそれぞれが、前記周辺装置を共用する、請求項1に記載のシステム。

【請求項4】前記接続が、USB仕様に従う、請求項1に記載のシステム。

【請求項5】前記システムが、前記システム・ノードと前記周辺装置とが接続されている時に、第2システム・ノードと前記周辺装置との間の接続を防止するステップを実行するようにも構成される、請求項1に記載のシステム。

【請求項6】前記装置コントローラが、複数のコンピュータ・システムに接続され、前記周辺装置が、前記複数のコンピュータ・システムの間で共用される、請求項1に記載のシステム。

【請求項7】前記コンピュータ・システムが、ラックマウント・システムである、請求項1に記載のシステム。

【請求項8】複数の周辺装置が接続される時に、前記周辺装置のすべてが、前記確立ステップ中に前記システム・ノードに接続される、請求項1に記載のシステム。

【請求項9】各仮想コンピュータ・システムが、少なくとも1つのシステム・プロセッサと、前記システム・プロセッサによって読み書きされるために接続されるメモリとを有する、複数の仮想コンピュータ・システムと、前記複数の仮想コンピュータ・システムに接続された、共通入出力コントローラと、

前記共通入出力コントローラに接続された、少なくとも1つの周辺装置とを含み、前記共通入出力コントローラが、前記複数の仮想コンピュータ・システムの1つによって要求された時に、前記仮想コンピュータ・システムと前記周辺装置との間の排他的接続を確立するマルチプロセッサ・コンピュータ・システム。

【請求項10】前記仮想コンピュータ・システムが、対称マルチプロセッサ・システム内で区分される、請求項9に記載のシステム。

【請求項11】前記共通入出力コントローラが、前記複数の仮想コンピュータ・システムに前記周辺装置への接続を提供するように接続される、請求項9に記載のシステム。

【請求項12】前記仮想コンピュータ・システムのそれぞれが、USBを介して前記共通入出力コントローラに接続され、各周辺装置が、USBハブを介して前記共通入出力コントローラに接続される、請求項9に記載のシステム。

【請求項13】前記仮想コンピュータ・システムのそれぞれが、USBコントローラを含む、請求項9に記載のシステム。

【請求項14】複数の周辺装置が接続される時に、前記複数の周辺装置のすべてが、前記周辺装置が接続される時に前記仮想コンピュータ・システムに接続される、請求項9に記載のシステム。

【請求項15】通信コントローラと、それぞれが機能的に前記通信コントローラに接続される、複数のノード入力と、

それぞれが機能的に前記通信コントローラに接続される、複数の装置接続とを含み、前記通信コントローラが、前記装置接続のノード所有権を調停し、所与のノード入力が前記装置接続の所有権を与えられる時に、そのノード入力と前記装置接続のすべてとの間の排他的接続を提供する装置コントローラ。

【請求項16】さらに、前記装置接続に接続された統合USBハブを含む、請求項15に記載の装置コントローラ。

【請求項17】前記装置接続が、USB接続である、請求項15に記載の装置コントローラ。

【請求項18】さらに、機能的に前記通信コントローラに接続されたISA装置接続を含む、請求項15に記載の装置コントローラ。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、全般的にはマルチプロセッサ・コンピュータ・システムに関し、具体的には、区分けされたマルチプロセッサ・システム内のプロセッサ間でのリソース割振りに関する。さらに具体的に言うと、本発明は、マルチプロセッサ・コンピュータ・システムのノード間で入出力装置の共通の組を共有するための方法に関する。

【0002】

【従来の技術】マルチプロセッサ・コンピュータ・システムは、当技術分野で周知であり、処理タスクを複数の異なるシステム・プロセッサ間で分割できるようにすることによって、処理能力の向上をもたらす。従来のシス

テムでは、各プロセッサが、システム・リソースのすべてにアクセスできる、すなわち、メモリ および入出力装置などのシステム・リソースのすべてが、システム・プロセッサのすべての間で共用される。通常、システム・リソースの一部を、プロセッサ間で区分することができ、たとえば、各プロセッサは共用メモリ にアクセスできるが、このメモリ が分割され、各プロセッサがそれ自体のワークスペースを有するようになっている。非一様メモリ・アクセス(Non-Uniform Memory Access : NU MA) システムでは、各プロセッサが、それ自体のメモリ を有し、他のプロセッサによって所有されるメモリ にもアクセスすることができる。

【 0003 】より最近になって、対称マルチプロセッサ(SMP) システムが、複数の独立コンピュータ・システムとして振る舞うように区分けされた。たとえば、8 プロセッサを有する単一のシステムを、処理のための別々のシステムとして8 つのプロセッサ(または1 つまたは複数のプロセッサの複数のグループ) のそれぞれを扱うように構成することができる。これらの「仮想」システムのそれぞれが、オペレーティング・システムのそれ自体のコピーを有し、独立にタスクを割り当てられるか、高速処理および改善された信頼性の両方を提供する処理クラスタとして一緒に動作することができる。通常、マルチプロセッサ・システムには、システム構成および特定のプロセッサと共用バスおよび共用装置との間のデータ・ルーティングを含む、システム全体の始動および動作を管理する、「サービス」プロセッサも存在する。

【 0004 】単一のマルチプロセッサ・システム内で、複数の仮想システムが、クラスタとして動作するように構成される時には、ソフトウェア・サポートを提供して、各クラスタ・ノードがマルチプロセッサ・システム内の他のノードのそれぞれと通信して、定足数(quorum) の協議および定足数の検証を実行し、拍動(heartbeat) 信号を送り、何らかのクラスタ通信技法を使用して他の定足数の機能を実行できるようにしなければならない。これが達成される時、もしプロセッサの1 つが障害を発生し、そのノードがクラスタから使用不要になるとして、そのノードに割り当てられたジョブを、標準的なクラスタ技法を使用して、残りのプロセッサ(ノード) の間で再割り当てすることができる。

【 0005 】典型的には、マルチプロセッサ・システムが、複数の仮想システムに分割される時には、仮想システムのそれぞれが、オペレーティング・システムのそれ自体のコピーを有し、同一のオペレーティング・システムが、仮想システムのそれぞれに使用される。各プロセッサは、同一のオペレーティング・システムを実行しているため、プロセッサの間でのリソース割振りを提供することは、比較的容易であった。

【 0006 】大規模マルチプロセッサ・システムの1 つ

の特徴として、それらが通常は大規模処理ジョブに使用されるので、キーボード、ディスプレイ、取外し可能媒体ドライブなどの通常の入出力装置を相対的にほとんど使用しない、ということがある。しかし、これらの装置は、頻度は低いが必要な時がありえるので、取り外すことができない。大規模マルチプロセッサ・システム内のノードのそれぞれでこれらの装置を可用とすることは、これらのほとんど使用されない装置の費用のかかる重複をもたらす、機器の管理および保守の不必要な負荷につながる。したがって、大規模マルチプロセッサ・システムの複数の区画またはノードが、単一の組の入出力装置を共用する手段が要望される。

【 0007 】

【 発明が解決しようとする課題】本発明の1 つの目的は、マルチプロセッサ・コンピュータ・システムを動作させるためのシステムおよび方法を提供することである。

【 0008 】本発明のもう1 つの目的は、マルチプロセッサ・コンピュータ・システム内の改善されたリソース割振りのシステムおよび方法を提供することである。

【 0009 】本発明のさらなる目的は、マルチプロセッサ・コンピュータ・システムのノード間で入出力装置の共通の組を共用するシステムおよび方法を提供することである。

【 0010 】

【 課題を解決するための手段】このように、本発明によれば、マルチプロセッサ・システムの複数のノードが1 組の入出力装置を共用できるようにするシステムおよび方法が提供される。すなわち、マルチプロセッサ・システム・ノードと共通の入出力装置の間の通信を管理し、個々のノードが1 つまたは複数のそのターゲット装置に排他的にアクセスできるようにする、キャビネット入出力コントローラ(CI / OC) が設けられる。ノードのそれぞれは、サービス・プロセッサを介してCI / OC と通信し、CI / OC は、さまざまな入出力装置とノードのUSB (universal serial bus) コントローラを相互接続する。別の実施形態では、レガシ入出力装置の接続を可能にするために、USB - I SAブリッジも含まれる。

【 0011 】

【 発明の実施の形態】好ましい実施形態は、マルチプロセッサ・システムの複数のノードが共通の入出力装置を共用できるようにするキャビネット入出力コントローラ(CI / OC) を提供する。好ましい実施形態では、イーサネット(登録商標) ・アダプタおよび関連するネットワーク接続などの他の実行時性能に関する入出力が、各ノードに存在し続ける。

【 0012 】ここで図面、具体的には図1 を参照すると、本発明の好ましい実施形態を実施することができるデータ処理システムのブロック図が示されている。デー

タ処理システム100は、たとえば、米国ニューヨーク州アーモンクのInternational Business Machines Corporation社から入手可能なコンピュータのサーバ・モデルの1つとすることができる。データ処理システム100には、プロセッサ101および102が含まれ、例示的实施形態では、プロセッサ101および102のそれぞれが、レベル2(L2)キャッシュ103および104に接続され、L2キャッシュ103および104は、システム・バス106に接続される。

【0013】システム・バス106には、システム・メモリ108およびプライマリ・ホスト・ブリッジ(PHB)122も接続される。PHB122は、入出力バス112をシステム・バス106に接続し、一方のバスから他方のバスへのデータ・トランザクションを中継および/または変換する。例示的实施形態では、データ処理システム100に、入出力バス112に接続され、ディスプレイ120用のユーザ・インターフェース情報を受け取るグラフィックス・アダプタ118が含まれる。ハード・ディスク・ドライブとすることのできる不揮発性ストレージ114、通常のマウス、トラックボール、または類似物とすることのできるキーボード/ポインティング・デバイス116などの周辺装置が、ISA(Industry Standard Architecture)ブリッジ121を介して入出力バス112に接続される。PHB122は、入出力バス112を介して、PCIスロット124およびUSB(Universal Serial Bus)コントローラ126にも接続される。

【0014】図1に示された例示的实施形態は、本発明の説明のためにのみ提供されるものであり、本発明は、図1に示されるような実施形態に限定されるものではない。たとえば、データ処理システム100に、コンパクト・ディスク読取専用メモリ(CD-ROM)ドライブまたはDVD(digital video disk)ドライブ、サウンド・カードおよびオーディオ・スピーカ、および多数の他の装置もオプションとして含めることができる。そのような変形形態のすべてが、本発明の趣旨および範囲に含まれると考えられる。下記のデータ処理システム100およびCI/OCアーキテクチャの例は、説明のための例としてのみ提供されるものであって、アーキテクチャの制限を意図する目的のものではない。

【0015】ここで図2を参照すると、大規模マルチプロセッサ・システムの基本構成要素として使用することができるノード200が示されている。ノード200には、SMPプロセッサ202および関連するメモリ204(他のプロセッサと共用することができる)が含まれる。SMPプロセッサ202は、プライマリ・ホスト・ブリッジ(PHB)206に接続され、PHB206は、USB(Universal Serial Bus)コントローラ208およびPCIスロット210に接続される。PCIスロット210には、この実施形態では他のノードと共用

されない、入出力装置212が接続(または、通常はプラグイン)される。USBコントローラ208は、CI/OC216の好ましいノード入力接続であり、CI/OC216は、サービス・プロセッサ(SP)214によって制御される。CI/OC216は、ノード入力218を介して他のノードに接続され、それらのノードが入出力装置220を共用できるようにする。この図には、1つの例示的ノードの詳細だけが示されているが、ノード200および他のノードが、CI/OC216の同一のノード入力によって、それぞれCI/OC216に接続されることに留意されたい。

【0016】大規模マルチプロセッサ・システムは、複数の小さい独立の区画として配置構成するか、NUMAまたはクラスタとして配置構成することができる。大規模マルチプロセッサ・システムが、NUMAまたはクラスタとして配置構成される場合、オプションのノード相互接続ハードウェア222を使用して、所望の相互接続構成を得ることができる。好ましい実施形態のCI/OC216は、ノードを制御し、装置の単一の共通の集合に相互接続するのに使用される。たとえば、単一のラックに収納された大規模マルチプロセッサ・システムには、複数のコンピュータ・ノードが含まれるが、そのラック全体が、単一のオペレータ端末、ディスク・ドライブなどだけを必要とする。SP214が、CI/OC構成を管理することが好ましい。各ノードには、USBコントローラ208が含まれる。

【0017】ここで図3を参照すると、好ましい実施形態によるCI/OCの高水準図が示されている。この図では、CI/OC300は、そのサービス・プロセッサ(SP)302に接続された状態で図示されている。例示的なノード入力306および308が図示されており、これらによって、図2に示されたものなどのノードの接続が可能になる。CI/OC300は、共通入出力装置304にも接続される。大規模マルチプロセッサ・システムごとに1つのCI/OC300と1つのSP302だけがあり、ノードのすべてが、CI/OC300の下流に接続される共通入出力装置304を共用することが好ましい。CI/OC300のノード入力数は、実装依存である。

【0018】ここで図4を参照すると、ノード入力406を下流の入出力装置404に接続するための一連のスイッチ408を含むものとして、CI/OC400のより詳細な図が示されている。SP402は、最大で1つのノードを、共通の入出力装置404の接続に切り替える。好ましい実施形態では、すべての接続がUSB準拠であることが規定されるので、装置は、ホット・プラグ可能すなわち、ノードの動作中にノードに接続し、取り外すことができる。スイッチ408の活動化は、下流接続の接続と同等である。スイッチの非活動化は、下流接続の取外しと同等である。SP402は、下流の入出力

装置404を、必要な時にノード入力406に接続する。

【0019】ここで図5を参照すると、共通入出力接続性を使用可能にするのに使用される基本入出力を含む共通入出力システムならびにさまざまな例示的入出力装置が示されている。この図では、CI / OC 500が、ハブ502に接続され、ハブ502は、好ましい実施形態ではUSBハブである。ハブ502は、CI / OCからの通信を、接続された装置に渡せるようにする。これらの装置には、キーボード504およびマウス506を含めることができる。要件を満たすネイティブUSB装置が入手可能でない場合には、入出力装置に、市販のUSB-I SA変換論理508を含めることができる。そうすることによって、I SAフロッピー・ドライブ510、シリアル・ポート512、およびパラレル・ポート514などのレガシ・デバイスを接続することができる。

【0020】<http://www.usb.org>で入手可能(本願の出願日時点)であり、参照によって本明細書に組み込まれるUniversal Serial Bus仕様では、VBUS、GND、D+、およびD-からなる4本のワイヤを介するUSB転送信号および電力が指定されている。シグナルは、2本のワイヤ、D+およびD-を介して行われる。好ましい実施形態では、ハブ502は、接続された装置に電力を供給する電力供給ハブ(powered hub)であり、その結果、ノードのUSBコントローラからのVBUSおよびGNDは必要ない。この実施形態では、図2に示されたノードのUSBコントローラ208からCI / OC 216への、および図5に示されたその後のハブ502への相互接続は、USBシグナルのD+およびD-からなり、VBUSおよびGNDは省略される。

【0021】ここで図6を参照すると、共通入出力システムの好ましい動作の流れ図が示されている。ノードは、共通入出力装置を必要とする時に(ステップ600)、サービス・プロセッサに要求を送って(ステップ610)、接続動作を要求する。他のノードが現在入出力チャネルを使用していない場合には(ステップ620)、SPは、CI / OCを切り替えて、そのノードが下流装置と通信できるようにする(ステップ630)。ノードは、入出力装置を使用し(ステップ640)、終了した時に、SPにCI / OC接続をオフに切り替えるように指示し、入出力装置を切り離す(ステップ650)。ノードは、通常の動作を継続する(ステップ660)。

【0022】入出力チャネルが別のノードによって使用中である場合には(ステップ620)、接続が拒否される(ステップ670)。ノードは、通常動作を再開し、必要な回数だけ接続の確立を再試行することができる。マルチプロセッサ・システムの性質のゆえに、この種の装置衝突は、比較的めずらしく、より洗練されたアービ

トリビション技法を実施することができる。

【0023】修正形態および変更形態好ましい実施形態に関して本発明を具体的に図示し、説明してきたが、形態および詳細におけるさまざまな変更を、本発明の趣旨および範囲から逸脱せずに行うことができることを、当業者は理解するであろう。たとえば、大量のノード数をサポートしなければならない場合に、CI / OCブロックをカスケード接続することができ、その結果、各ノードが1つのCI / OCに接続され、CI / OC装置のチェーンを介して周辺装置と通信できるようにすることができる。上記および他の修正形態は、請求の範囲内であるものとみなされる。

【0024】まとめとして、本発明の構成に関して以下の事項を開示する。

【0025】(1)それぞれが少なくとも1つのシステム・プロセッサおよび関連メモリを有し、それぞれがそれぞれの装置ポートを介して通信するように接続された、複数の装置ノードと、前記装置ポートの少なくとも1つに接続された、装置コントローラと、前記装置コントローラに接続された、少なくとも1つの周辺装置とを含み、前記コンピュータ・システムが、直接接続を介してシステム・ノードから前記装置コントローラへ要求を送るステップと、前記装置コントローラによって、前記システム・ノードと前記周辺装置との間の排他的接続を確立するステップと、前記システム・ノードによって前記周辺装置を操作するステップとを実行するように構成されるコンピュータ・システム。

(2)前記システム・ノードが、対称マルチプロセッサ・システムとして動作する、上記(1)に記載のシステム。

(3)前記システム・ノードのそれぞれが、前記周辺装置を共用する、上記(1)に記載のシステム。

(4)前記接続が、Universal Serial Bus仕様に従う、上記(1)に記載のシステム。

(5)前記システムが、前記システム・ノードと前記周辺装置とが接続されている時に、第2システム・ノードと前記周辺装置との間の接続を防止するステップを実行するようにも構成される、上記(1)に記載のシステム。

(6)前記装置コントローラが、複数のコンピュータ・システムに接続され、前記周辺装置が、前記複数のコンピュータ・システムの間で共用される、上記(1)に記載のシステム。

(7)前記コンピュータ・システムが、ラックマウント・システムである、上記(1)に記載のシステム。

(8)複数の周辺装置が接続される時に、前記周辺装置のすべてが、前記確立ステップ中に前記システム・ノードに接続される、上記(1)に記載のシステム。

(9)各仮想コンピュータ・システムが、少なくとも1つのシステム・プロセッサと、前記システム・プロセ

10

20

30

40

50

サによって読み書きされるために接続されるメモリとを有する、複数の仮想コンピュータ・システムと、前記複数の仮想コンピュータ・システムに接続された、共通入出力コントローラと、前記共通入出力コントローラに接続された、少なくとも1つの周辺装置とを含み、前記共通入出力コントローラが、前記複数の仮想コンピュータ・システムの1つによって要求された時に、前記仮想コンピュータ・システムと前記周辺装置との間の排他的接続を確立するマルチプロセッサ・コンピュータ・システム。

(10) 前記仮想コンピュータ・システムが、対称マルチプロセッサ・システム内で区分される、上記(9)に記載のシステム。

(11) 前記共通入出力コントローラが、前記複数の仮想コンピュータ・システムに前記周辺装置への接続を提供するように接続される、上記(9)に記載のシステム。

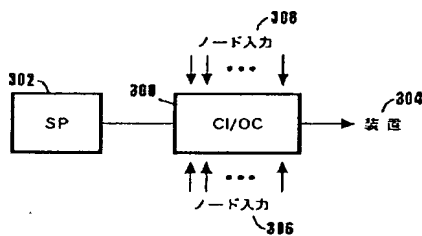
(12) 前記仮想コンピュータ・システムのそれぞれが、USBを介して前記共通入出力コントローラに接続され、各周辺装置が、USBハブを介して前記共通入出力コントローラに接続される、上記(9)に記載のシステム。

(13) 前記仮想コンピュータ・システムのそれぞれが、USBコントローラを含む、上記(9)に記載のシステム。

(14) 複数の周辺装置が接続される時に、前記複数の周辺装置のすべてが、前記周辺装置が接続される時に前記仮想コンピュータ・システムに接続される、上記(9)に記載のシステム。

(15) 通信コントローラと、それぞれが機能的に前記通信コントローラに接続される、複数のノード入力と、それぞれが機能的に前記通信コントローラに接続される、複数の装置接続とを含み、前記通信コントローラが、前記装置接続のノード所有権を調停し、所与のノード入力が前記装置接続の所有権を与えられる時に、そのノード入力と前記装置接続のすべてとの間の排他的接続

【 図3 】



を提供する装置コントローラ。

(16) さらに、前記装置接続に接続された統合USBハブを含む、上記(15)に記載の装置コントローラ。

(17) 前記装置接続が、USB接続である、上記(15)に記載の装置コントローラ。

(18) さらに、機能的に前記通信コントローラに接続されたISA装置接続を含む、上記(15)に記載の装置コントローラ。

【 図面の簡単な説明 】

【 図1 】 本発明の好ましい実施形態によるマルチプロセッサ・コンピュータ・システムを示す図である。

【 図2 】 本発明の好ましい実施形態に従ってCI/OCに接続されたSMPノードを示す図である。

【 図3 】 本発明の好ましい実施形態によるCI/OCのブロック図である。

【 図4 】 本発明の好ましい実施形態によるCI/OCおよびサービス・プロセッサのブロック図である。

【 図5 】 本発明の好ましい実施形態による、CI/OCの下流の、ハブに接続された複数の入出力装置のブロック図である。

【 図6 】 本発明の好ましい実施形態による、共通入出力システムの使用の流れ図である。

【 符号の説明 】

200 ノード

202 SMPプロセッサ

204 メモリ

206 プライマリ・ホスト・ブリッジ(PHB)

208 USB(Universal Serial Bus)コントローラ

210 PCIスロット

212 入出力装置

214 サービス・プロセッサ(SP)

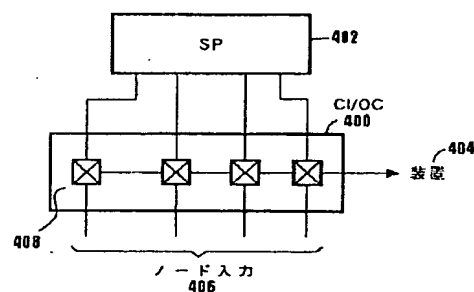
216 CI/OC

218 ノード入力

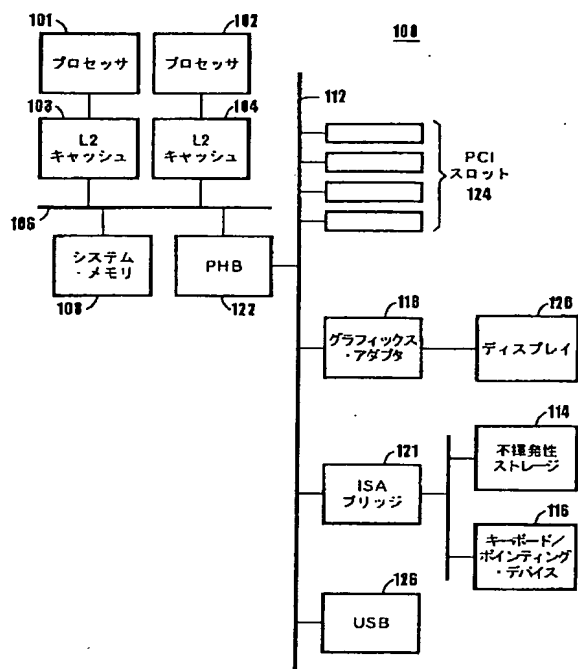
220 入出力装置

222 ノード相互接続ハードウェア

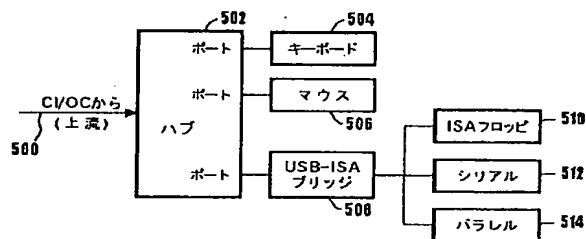
【 図4 】



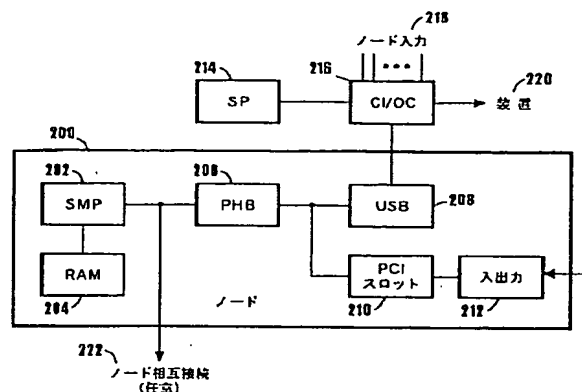
【 図1 】



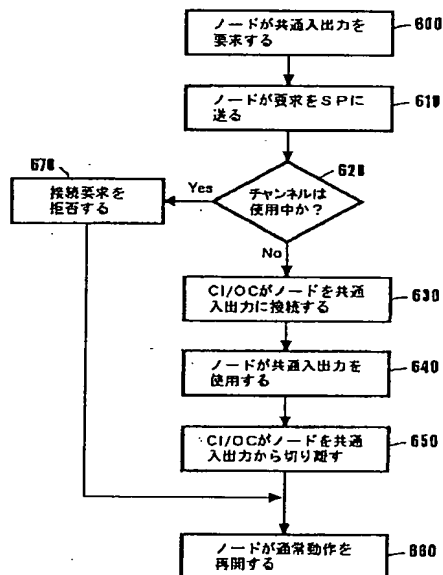
【 図5 】



【 図2 】



【 図6 】



フロント ページの続き

(72)発明者 リチャード・ビールコウスキー
 アメリカ合衆国98052 ワシントン州レッド
 モンド ワンハンドレッド・フィフティ・
 エイス・プレイス・エヌ・イー 8336

(72)発明者 パトリック・エム・ブランド
 アメリカ合衆国27613 ノースカロライナ
 州ローリー ウィロウド・コート 8904
 Fターム(参考) 5B014 HA07 HC13 HC15
 5B045 EE08 EE11